

1 **TITLE OF THE INVENTION**

2 **MEMORY SYSTEM FOR IMPROVING DATA INPUT/OUTPUT**

3 **PERFORMANCE AND METHOD OF CACHING**

4 **DATA RECOVERY INFORMATION**

5 a **OF**
6 **CLAIM FOR PRIORITY**

7 This application makes reference to, ~~incorporates the same herein~~, and claims all benefits
8 accruing under 35 U.S.C. §119 from an application for *MEMORY SYSTEM FOR IMPROVING*
9 *DATA INPUT/OUTPUT PERFORMANCE AND METHOD OF CACHING DATA RECOVERY*
10 *INFORMATION* earlier filed in the Korean Industrial Property Office on the 16th of September 1996,
11 and there duly assigned Serial No. 40202/1996, ~~a copy of which application is annexed hereto.~~

12 **BACKGROUND OF THE INVENTION**

13 Fiel/d of the Invention
14 Technical Field

15 The present invention relates to a memory system such as ~~redundant arrays~~ of inexpensive
16 disks (RAID) and, more particularly, to a redundant ~~array~~ of inexpensive disks capable of providing
17 high ~~performance of data input/output operation~~ and a method of caching data recovery information
18 using the redundant ~~array~~ of inexpensive disks.

19 Description of the Related Art
20 Related Art

21 A high

22 High technology computer system depends considerably on its central processor unit (CPU)

1 and input/output subsystem to increase overall system performance. While ^{the} information processing
2 speed of the CPU has been dramatically improved in recent years because of VLSI technology, the
3 performance of the input/output subsystem has not improved as desired. This increases the time
4 required to access data in the memory system. Furthermore, ~~as the cost needed to restore data is~~
5 increased when an error is generated in the input/output subsystem, an input/output subsystem
6 having excellent performance and reliability is needed. As a solution to this, a disk array system
7 known as a redundant array of inexpensive disks (RAID), constructed of a number of relatively small
8 capacity disk drives has been proposed as a low cost alternative to a single large expensive disk for
9 ~~storage~~
~~providing large storage of~~ digital information.

10
11 RAID systems are now commercially available as cost effective mass storage providing
12 reliable and continuous services to a host computer or network file server. The theory of RAID is
13 to use relatively inexpensive disks, which may individually have a higher chance of failure than
14 expensive disks, and compensating for this higher failure rate by adding redundancy by creating and
15 storing parity blocks to facilitate recovery from a disk failure. Reports on the performance and
16 reliability of disk arrays are presented in "*A Case For Redundant Arrays Of Inexpensive Disks*
17 (*RAID*)" by D. Patterson, G. Gibson, and R. H. Kartz, at Report No. UCB/CSD 87/89, December
18 1987, Computer Science Division (EECS), University of California, Berkeley, Calif. 94720.
19 Exemplars of contemporary RAID systems are disclosed in U.S. Patent No. 5,257,367 for *Data*
20 *Storage System With Asynchronous Host Operating System Communication Link* issued to
21 Goodlander et al., U.S. Patent Nos. 5,367,669 and 5,455,934 for *Fault Tolerant Hard Disk Array*

✓ GJP
✓ GJP

1 Controller issued to Holland et al., U.S. Patent No. 5,418,921 for *Method And Means For Fast*
2 *Writing Data To LRU Cached Based DASD Arrays Under Drivers Fault Tolerant Modes* issued to
3 Cortney et al., U.S. Patent No. 5,463,765 for *Disk Array System, Data Writing Method Thereof, And*
4 *Fault Recovering Method* issued Kakuta et al., U.S. Patent No. 5,485,598 for *Redundant Disk Array*
5 *(RAID) System Utilizing Separate Cache Memories For The Host System And The Check Data*
6 issued to Kashima et al., U.S. Patent No. 5,522,032 for *RAID Level 5 With Free Blocks Parity Cache*
7 issued to Franaszek et al., U.S. Patent No. 5,530,948 for *System And Method For Command Queuing*
8 *On RAID Levels 4 And 5 Parity Drives* issued to Islam, U.S. Patent No. 5,579,474 for *Disk Array*
9 *System And Its Control Method* issued to Kakuta et al., U.S. Patent No. 5,640,506 for *Integrity*
10 *Protection For Parity Calculation For RAID Parity Cache* issued to Duffy, and U.S. Patent No.
11 5,636,359 for *Performance Enhancement System And Method For A Hierarchical Data Cache Using*
12 *A RAID Parity Scheme* issued to Beardsley et al.

13 As generally discussed in the Patterson report and subsequent contemporary RAID systems
14 ~~as set forth~~, the large personal computer market has supported the development of inexpensive disk
15 drives having a better ratio of performance to cost than single large expensive disk systems. The
16 number of input/outputs (I/Os) per second per read/write head in an inexpensive disk is within a
17 factor of two of the large disks. Therefore, the parallel transfer from several inexpensive disks in
18 a RAID system, in which a set of inexpensive disks function as a single logical disk drive, produces
19 ~~better~~ greater performance than a single large expensive disk (SLED) at a reduced ~~price~~ ^{cost}.

1 Unfortunately, when data is stored on more than one disk, the mean time to failure varies
2 inversely with the number of disks in the array. In order to correct for this decreased mean time to
3 failure of the system, error recognition and correction is characteristic of all RAID systems.
4 ~~six structures~~
5 ~~Generally, each RAID system is organized in six structures commonly referred to as six levels each.~~
having a different means for error recognition and correction as described hereinbelow.

6 In RAID structure of level 0, data is distributed and stored in all drives in the disk array,
7 taking interests in performance rather than data reliability.

8 In RAID structure of level 1, the mirroring, a conventional method of improving the disk
9 ~~needs a lot of costs~~ ^{has a high cost} performance, needs a lot of costs since all contents of the disk must be stored in a reproduction disk
10 without change. Accordingly, in a database system requiring a large-capacity disk space, only the
11 ^{percent} fifty percent of the disk space can be used. However, the mirroring is the best way to enhance the
12 data reliability because identical data is stored in the reproduction disk. In RAID structure of level
13 2, this is used to minimize the cost required to enhance data reliability. The RAID structure of level
14 2 distributes and stores data in each disk array in ~~bits~~ ^{bytes}, and has several test disks using a ~~hamming~~ ^{Hamming}
15 ~~code~~ code, besides the data disk, in order to recognize and correct errors.

16 In RAID structure of level 3, data is input/output in parallel to/from the drive when
17 input/output is requested once, and parity data is stored in a separate drive. Furthermore, disk
18 spindles are synchronized so as to make all drives simultaneously input or output data. Accordingly,

1 rapid data transmission can be carried out even if parallel input/output is not performed fast. If one
2 drive has an error, the erroneous data can be restored by using the currently operated drive and parity
3 drive even though the total data rate is decreased. The RAID structure of level 3 is used in an
4 ~~application which requires very fast data transmission rate, super computer and image manipulation~~
5 processors. That is, the RAID of level 3 has higher efficiency in a long data block transmission but
6 has lower efficiency in a short data block transmission which requires fast input/output request.
7 Furthermore, since the data drive is used together with a single drive for redundancy, ~~the drive~~
8 smaller than that used in the RAID of level 1 is used but its controller becomes more expensive and
9 complicated.

10 In RAID structure of level 4, the parity data is calculated and stored in a separate drive, and
11 data is striped across. The data can be restored when it has error. Its reading performance is similar
12 to that of RAID of level 1 but its writing is much poorer than the single drive because the parity
13 information must be provided to the single drive. Thus, the RAID structure of level 5 having
14 improved writing performance is supplemented to the RAID of level 4.

15 In RAID structure of level 5, data is striped across in each drive array, and parity data is
16 distributed and stored in all drives in order to remove ~~bottleneck~~ bottleneck phenomenon when data is written.
17 In this RAID structure, since the data written in all drives must be read in order to calculate the parity
18 when the data is written, its speed is slower. However, it is possible to process the data input/output
19 transmission and to restore data stored in a drive having an error. Accordingly, the RAID structure of

at
BB 1 level 5 is effective in recording of long data, and also effective in recording of short data if an application program gives weight *to* the data reading or the array design is improved in order to increase the writing performance. Even if the size of the data block is decreased, performance and data availability can be obtained to some degree. Moreover, the RAID structure of level 5 is most effective in terms of cost in comparison with a non-array device.

6 Among all disk array structures, the RAID structure of level 5 provides a higher reliability
7 with smaller additional cost, and at the same time, makes the parallel disk access possible, resulting
8 in the improvement of data processing rate. Generally, when data writing instruction is received
9 from the host computer for writing in each drive in the RAID structure of level 5, the CPU
10 determines a target location, and transmits the data to ~~controller~~ where old data and old parity
11 stored in each drive are read. ~~Controller~~ calculates a new parity based on an exclusive OR arithmetic
12 operation, and writes new data and new parity in a predetermined drive. However, when writing
13 instruction of a short data block is received from the host computer in the RAID structure of 5 level,
14 access of another disk on the strip is brought about which ~~leads~~ to a deterioration of the entire
15 system performance. I have observed that this phenomenon ~~remarkably~~ appears in the on-line
16 transaction processing environment having many operation loads. That is, in case of the partial strip
17 writing, old parity and old data are read from a predetermined drive, ~~an~~ exclusive-OR operation is
18 performed to determine new data, and then new parity information and new data are written in the
19 predetermined drive. Two-time reading and writing operations are necessarily required which results
20 in a larger overhead of data write in comparison with a single large expensive drive.

SUMMARY OF THE INVENTION

Accordingly, it is therefore an object of the present invention to provide a redundant arrays of inexpensive disks (RAID) system with an enhanced process performance and a reduced overhead of data write.

It is also an object to provide a RAID system capable of reducing an overhead during a read operation of data recovery information in order to improve its data input/output performance, and a method of caching data recovery information using the memory system.

These and other objects of the present invention can be achieved by a redundant arrays of inexpensive disks (RAID) system which includes a plurality of defect-adaptive memory devices for sequentially storing information needed for data recovery in a predetermined region of a recording medium in the form of block, and storing data in a region other than the predetermined region. A plurality of caches are connected to the adaptive memory devices to store information blocks needed for data recovery, the information blocks being read from a predetermined memory device. A controller is connected to each adaptive memory device and cache to control the writing and reading of data and information needed for data recovery in each memory device, calculate information needed for recovery of data read from each memory device, and store the information needed for recovery of data calculated in a predetermined cache.

The present invention is more specifically described in the following paragraphs by reference to the drawings attached only by way of example.

1 **BRIEF DESCRIPTION OF THE DRAWINGS**

2 A more complete appreciation of the present invention, and many of the attendant advantages
3 thereof, will become readily apparent as the same becomes better understood by reference to the
4 following detailed description when considered in conjunction with the accompanying drawings in
5 which like reference symbols indicate the same or similar components, wherein:

6 FIG. 1 is a block diagram of a RAID system;

7 FIG. 2 illustrates an example of data transmission of the RAID system;

8 FIG. 3 is a flowchart illustrating a process of writing data and parity information transmitted
9 from a host computer to each drive in the RAID system;

10 FIG. 4 is a block diagram of an RAID system constructed according to an embodiment of the
11 present invention; and

12 FIG. 5 is a flowchart illustrating a process of writing data and parity information in the RAID
13 system constructed according to the embodiment of the present invention.

14 **DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENT**

15 Referring now to the drawings and particularly to FIG. 1, which illustrates a redundant *array*
16 of inexpensive disks (RAID) system in level 5. As shown in FIG. 1, the RAID system includes a
17 central processing unit (CPU) 2, a controller 6 connected to the CPU 2 via an input/output bus 4, and
18 a plurality of disk drives DR1-DR5 connected to the controller 6 via SCSI bus 8.

19 CPU 2 transmits data transmitted through an input/output bus 4 from a host computer (not

1 shown) to the controller 6. The controller 6 connected to input/output bus 4 is controlled by CPU
2 to control input/output data between drive disks DR1 to DR5 which are connected to CPU 2 and
3 SCSI bus 8. Each drive DR1 to DR5 connected to SCSI bus 8 records and reproduces the data
4 transmitted from the host computer under the control of controller 6.

5

6 FIG. 2 illustrates an example of data transmission of the RAID structure in level 5. Data ND
7 transmitted from the host computer is divided by strip (the data is divided by strip 3 in FIG. 2),
8 distributed and stored in each drive DR1 to DR5. That is, each drive DR1 to DR5 has a data block
9 D in which data is stored, and a parity block P in which parity information is stored, to thereby store
10 the data transmitted from the host computer under the control of controller 6.

a flowchart

11 FIG. 3 is ~~a control flow chart~~ for explaining the writing of the data and parity information
12 transmitted from the host computer in each drive in the RAID structure of level 5. Referring to FIG.
13 3, when ^a data writing instruction is received from the host computer, ^{the} CPU 2 calculates a target
14 location at step 10. At step 12, CPU 2 transmits the data transmitted from the host computer to
15 controller 6. Controller 6 reads old data OD and old parity OP stored in each drive at steps 14 and
16 16. Next, the controller 6 calculates a new parity NP according to the following formula (1).

17
$$NP = OP \vee OD \vee ND \quad (\vee \text{ means exclusive OR}) \quad (1)$$

18 Controller 6 writes data ND and new parity NP in a predetermined drive at steps 20 and 22.

1 As described, ~~in case that~~ when a writing instruction of a short data block is received from the host
2 computer in the RAID system of level 5 structure, access of another disk on the strip is brought about
3 which ~~attributes~~ leads to a deterioration to the entire system performance. This ~~remarkably~~ appears in the
4 on-line transaction processing environment having many operation loads. That is, in case of the
5 partial strip writing, old parity OP and old data OD are read from a predetermined drive, exclusive-
6 ORed according to formula (1), its result is exclusive-ORed with data ND, and then new parity NP
7 and new data ND are written in a predetermined drive. Thus, two-time reading and writing
8 operations are needed which results in a larger overhead of write data in comparison with a single
9 large expensive drive.

10 Turning now to FIG. 4 which illustrates a RAID system to which parity cache arrays 38 are
11 connected according to an embodiment of the present invention. Referring to FIG. 4, the RAID
12 system consists of a CPU 30 for controlling the overall system. A controller 34 which is connected
13 to CPU 30 through an input/output bus 32 to distribute and store data transmitted from a host
14 computer to each drive array 39, or reproduce the stored data under the control of ~~CPU 30~~. Drives
15 1 to 5 (39) which are connected to controller 34 through SCSI bus 36 to store and reproduce the data
16 and data recovery information (parity information) transmitted from the host computer under the
17 control of controller 34. Caches 1 to 5 (38) which are connected to controller 34 and input/output
18 bus 36 placed between drives 39 to store the parity information.

19 Each drive 39 consists of a plurality of blocks in order to store and read the data and parity

1 information. Furthermore, each drive 39 sets up the predetermined number of parity block from the
2 cylinder zero on the disk, and uses it as a parity information storing region, without using the
3 stripping method defined in the RAID structure in level 5. Here, the data cannot be recorded in the
4 parity information storing region.

a flowchart

5 FIG. 5, a control flow chart for explaining a process of writing data and parity information
6 in the RAID system constructed according to the embodiment of the present invention. The control
7 process of writing data will be explained in detail with reference to FIGS. 4 and 5 hereinbelow.

8 First of all, the data writing instruction is received from the host computer, CPU 30 updates
9 a task file required at step 40, and then calculates a target cylinder (=parity block + request cylinder)
10 in order to use a separate parity block in the drive. Then, CPU 30 transmits new data ND to be
11 written at step 42. Controller 34 next reads old data OD from a predetermined drive 39 in order to
12 generate new parity NP, and then examines if old parity information OP to be read is hit in cache 38
13 at step 46. Here, if the old parity information OP is hit in cache 38, controller 34 proceeds to step
14 50. If the old parity information is not hit in cache 38, controller 34 proceeds to step 48. That is,
15 in case that the old parity information OP and parity information are not hit, controller 34 reads the
16 old parity information OP from the predetermined drive 39 at step 48, updates a cache table, and then
17 moves to step 50. Controller 34 calculates a new parity NP by exclusive-ORing the old parity
18 information read and the new data ND through the following formula (2).

1 NP=OP\OD\ND ----- (2)

The controller or loads and predetermined cache 38
2 Controller 34 updates the cache table at step 52, and then writes the new data ND transmitted
3 from the host computer and the calculated new parity NP in a predetermined drive at steps 54 and
4 56. Then, the data writing process of the present invention is completed.

5 According to the present invention, the parity cache is connected between each drive and
6 ~~controller in order to rapidly apply parity information read request. Furthermore, since the parity~~
7 ~~block for storing the parity information is set up from the cylinder zero on the disk, it is now possible~~
8 ~~to prevent time delay due to a separate search when sequential read/write operation is carried out.~~

9 While there have been illustrated and described what are considered to be preferred
10 embodiments of the present invention, it will be understood by those skilled in the art that various
11 changes and modifications may be made, and equivalents may be substituted for elements thereof
12 without departing from the true scope of the present invention. In addition, many modifications may
13 be made to adapt a particular situation to the teaching of the present invention without departing
14 from the central scope thereof. Therefore, it is intended that the present invention not be limited to
15 the particular embodiment disclosed as the best mode contemplated for carrying out the present
16 invention, but that the present invention includes all embodiments falling within the scope of the
17 appended claims.